

# The Princeton Lexicon Project: A Report on WordNet

George A. Miller, Christiane Fellbaum, Judy Kegl, and Katherine Miller

## Introductory Note

This research was supported in part by contract N00014-86-K-0492 between the Office of Naval Research and Princeton University, in part by contract MDA 903-86-K-0242 between the Army Research Institute and Princeton University, and in part by a grant from the James S. McDonnell Foundation to Princeton University. The views and conclusions contained herein are those of the authors and should not be interpreted as representing official policies, either expressed or implied, of the Office of Naval Research, the Army Research Institute, the McDonnell Foundation, or Princeton University. WordNet represents research in progress and is not available for distribution. WordNet has benefitted from the input of numerous lexicographers and programmers in addition to the authors: Amalia Bachman, Marie Bienkowski, Richard Beckwith, George Collier, Melanie Cook, Derek Gross, Dan Teibel, and Anton Vishio.

## Introduction

WordNet is an electronic lexical reference system for English, designed in accordance with psycholinguistic theories of the organization of human lexical memory. This novel lexical reference system for English is being developed in the form of an electronic database. Its design derives from psychological and linguistic theories about how lexical information is organized and stored in the memories of people who know English well and speak it fluently. The success of this experimental system would demonstrate the adequacy of the theories from which it derives, but even if those theories must be revised or replaced, the lexical database that is being developed in order to test them will be adaptable to a variety of practical applications. WordNet, supplemented on-line by machine-readable dictionaries and made available via a multi-window workstation, can be profitably incorporated into any task that is facilitated by easy access to lexical information.

Word knowledge is analyzed into: 1) the sound pattern, 2) the concept that the sound pattern can express, and 3) the association of sound and concept. Sounds and concepts are learned differently: as a consequence, different kinds of lexical relations are established: 1) phonological (e.g., rhyme) and morphological relations (e.g. inflection, derivation, compounding) are word-specific, whereas 2) semantic relations (e.g., synonymy, subordination, part-whole) are truth-functional.

Both kinds of relations are incorporated in WordNet. A concept is represented by a set of synonyms that can be used, in appropriate contexts, to express it; other semantic relations are represented by labeled pointers between the related concepts. WordNet will test the adequacy of current ideas about the structure of the lexicon

by testing whether a realistically large sample of the English lexicon can be represented in this way.

The use of synonym sets is both an innovative and an expedient approach to dictionary design. Standard dictionaries develop uniform semantic representations for all the lexical items in English by systematizing the writing of sense definitions or by determining a set of linguistic primitives that constitute the meaning of lexical items. WordNet circumvents the writing and systematizing of sense definitions by representing concepts as relations among words arranged in a "vocabulary matrix," a giant network coding various relations by means of connections between words. It simply looks along a given row of the vocabulary matrix, notes all the words that can be used to express the same concept, and then substitutes that synonym set for the statement of the concept. If one accesses the dictionary by way of the horizontal word list, one gets a view of the polysemy of a word (all the different concepts that the word can be associated with). On the other hand, if one accesses the matrix from the vertical concept list, one gets a row containing all the different synonymous words that express a given concept.

Once the basic matrix is in place, an elaborate system of cross-referencing allows the coding of various relations between synonym sets, including relations of antonymy, superordination, subordination, part-whole, grading, and presupposition. WordNet is free from any requirement to encode all the information about a word in the confines of a single entry. Furthermore, the nonlinear nature of this net together with the freedom afforded by computer access captures many important relations obscured by the formatting constraints of hand-held dictionaries.

### **Psycholinguistic Issues**

What a language user must know and how that knowledge is organized are related but separable questions. In order to speak and understand any language, it is necessary to know the sounds and meanings of thousands of different lexical units.

How that lexical knowledge is organized, however, is a much more difficult question. Whereas in a printed dictionary it is organized alphabetically, in a person's memory the organization is much more complex. Lexical memory must be so organized that the sounds and the contextually appropriate meanings of thousands of different words can be retrieved at rapid rates, otherwise the conversational use of language would scarcely be possible. The nature of this organization and how it comes to be constructed during the process of learning a language are basic questions for psycholinguistic research. Questions about the organization of lexical memory are easier to consider, however, if one first becomes clear about what a language user must know.

A vocabulary matrix is sufficiently general to represent any lexicon, whether it exists in a person, in a book, or in a computer. It contains a representation of the phonological form of a word and a representation of the conceptual content of the word, along with the associative bond connecting them. The vocabulary matrix is not a complete model of a human language user's lexical knowledge, however. A good model of a person's lexical knowledge would have to include the phonological and morphological features of the words and the semantic and pragmatic relations among lexical concepts.

### *Lexical Relations*

Not only are the phonological and morphological relations that exist between words not shown in the vocabulary matrix, but conceptual relations are not represented, either. A wide variety of such relations have been studied by psycholinguists (Chaffin & Herrmann 1984). For example, subordination and superordination (e.g., a maple is a tree, and a tree is a plant), which are relations between concepts, do not appear in a simple listing of lexical concepts. These relations, referred to as *hyponymy*, generate a hierarchical structure, a taxonomy, in the lexicon.

The part-whole relation, or *meronymy*, is also a relation between concepts, not between words (Iris, Litowitz, & Evens 1985). For example, a car has an engine, an engine has a carburetor, and a carburetor has a flutter valve; that is, *flutter valve* is a meronym of *carburetor* and *carburetor* is a meronym of *engine*. Like hyponymy, meronymy exhibits a hierarchical organization where, instead of the ISA relation, the HASA relation is exploited.

No adequate theory of the organization of lexical memory can ignore the strong formal relations between the columns or the strong semantic relations between the rows of the vocabulary matrix. Lexical relations must, therefore, be included in any electronic system that hopes to simulate the structure of human memory. The vocabulary matrix is merely a skeleton; it must be fleshed out with many formal and conceptual relations.

### *Sources of Evidence*

Two rather different kinds of factual data are available to support claims about the organization of lexical memory. One is linguistic: the data underlying theories of lexical organization are conveniently summarized in printed dictionaries and thesauruses. The second is psychological: a variety of experimental investigations have provided evidence for the psychological reality of the hypothesized mental structures.

The linguistic evidence comes primarily from dictionaries and thesauruses that summarize the relevant linguistic information derived ultimately from the recorded use of the language by native speakers and from native speakers' subjective judgments.

In addition to corpus-based lexicography, some linguists, lexicographers and psychologists also rely on native speaker intuitions. For example, psychologists sometimes give native speakers a word and ask what other words it suggests, or they may constrain the person's associations by specific instructions, such as "What is a kind of plant?" or "List all the trees you can think of." Judgments that ISA or HASA relations hold take the form of judgments of the truth or falsity of such statements as "A maple is a tree" or "A gasoline engine has a carburetor." Linguists are more likely to frame questions in terms of sentences, such as "Do  $S_1$  and  $S_2$  have the same meaning?" where  $S_1$  and  $S_2$  are identical sentences except for a pair of words whose meanings are to be compared. Or they may ask for judgments of oddness, for example, "pines and other maples" sounds odd, "trees and other maples" sounds odd, but "pines and other trees" does not.

The conceptual dimension of lexical memory has also been explored experimentally by psychologists. One of the landmark studies was the work of Collins and Quillian (1969) who reported that it takes people longer to judge the truth of the statement *A canary is an animal* than to judge *A canary is a bird*. They attributed such observations to the fact that *bird* is the immediate superordinate of *canary*, whereas *animal* is a more remote superordinate.

Although the work outlined in this paper is not basic research in the sense that the experimental studies just mentioned clearly are, it can nevertheless contribute to the understanding of the organization of lexical memory. The contribution follows from the inclusion of a sizeable fraction of the English lexicon. Psychological experiments are almost necessarily conducted with a small number of words and then assumed (often implicitly) to generalize over the entire vocabulary. A failure to look for negative evidence can tempt one into an (unconscious) favoritism for words that confirm one's hypothesis. Therefore, it is advisable to test hypotheses against a large collection of words, a collection assembled in ignorance of the hypotheses in question.

### **WordNet: Implementation of a Model of Lexical Organization**

WordNet is an electronic lexical reference system designed in accordance with the theories summarized above. The first step in creating WordNet was to invent an electronic version of the vocabulary matrix.

#### *Synonym Sets*

A major problem in constructing a vocabulary matrix is how to represent all the various concepts that words can express.

Lexicographers represent lexical concepts by circumlocution, that is to say, they use words to define words. Lexicographers take great pains to distinguish among different senses that a given word can express, but they pay far less attention to establishing a common phrasing for the same sense when it appears in entries for different words. For example, in one widely used dictionary the same lexical concept is phrased as "inferior in quality or value" in the definition of *poor* and as "of little or less importance, value, or merit," in the definition of *inferior*. If WordNet represented the lexical concepts in the vocabulary matrix by definitional phrases borrowed from a conventional dictionary, many, perhaps most, synonymic relations would be overlooked.

Some standard convention for expressing word senses is required. Many different notations for lexical concepts have been proposed (see, for example, Anderson 1976; Jackendoff 1983; Katz 1972; Miller & Johnson-Laird 1976; Norman & Rumelhart 1975; Schank 1972; Sowa 1984; Talmy 1985), but they have been worked out in detail for only small sets of English words.

In order to proceed with WordNet, we have used *synonym sets* to represent lexical concepts. That is to say, the identifier for the concept on any given row of the vocabulary matrix is given by the list of words that (in appropriate contexts) can be used to express that concept. Since the synonym sets will be numbered, each concept

will be represented in the system by a number, but displayed to the user as a set of words having a shared meaning.

It should be noted that synonym sets, unlike dictionary entries, do not have headwords. In a book of synonyms, for example, one entry might have *pipe* as the headword, alphabetized under *P* with “tube” as its contents, and another entry might have *tube* as the headword, alphabetized under *T* with “pipe” as its contents. In WordNet, the synonym set {pipe, tube,} stands as an elementary component, and neither word is ahead of the other. This practice has the advantage of symmetry: if *x* is a synonym of *y*, then *y* is a synonym of *x*.

Because synonymy is so central to the design of WordNet, it resembles the electronic thesauruses that are now becoming available commercially (Raskin 1987). WordNet goes beyond those products, however, by incorporating conceptual relations other than synonymy — as will be described.

### *The Master List*

Once a satisfactory list of synonym sets becomes available, it will be simple to index it. That is to say, an alphabetical listing of all the words in all the synonym sets can be constructed where each word is followed by the numbers of all the synonym sets of which it is a member. This “master” list, can also contain information that is word-specific and not dependent on the concepts that the word can be used to express, such as the relative frequency of use of each word.

### *Conceptual Relations*

As of the end of 1987, the WordNet files included 11,500 different nouns organized into over 7,000 synonym sets; approximately 6,000 different verbs organized into over 3,000 synonym sets; and 9,500 different adjectives organized into over 8,200 synonym sets. That gave a total of over 27,000 different words organized into approximately 18,200 synonym sets. The next step was to introduce both semantic relations between lexical concepts (Cruse 1986; Evens et al. 1983; Lyons 1977, ch. 9) and other relations as well. While additional synonym sets continue to be added, we are now introducing cross-references designed to represent conceptual relations.

Conceptual relations are represented in WordNet by cross-references between synonym sets. Each synonym set, therefore, will be followed by a list of the numbers of other synonym sets related to it in particular ways.

Hyponymy, for example, can be introduced in WordNet by appending to a given synonym set one number that points to its superordinate term and other numbers that point to its hyponyms. Meronymy, too, is introduced in WordNet by labeled cross-references.

### *The Hyponymic Hierarchy*

Cognitive psychologists have been interested in lexical hierarchies at least since Collins and Quillian (1969) proposed them as a model of semantic memory. According to the theory, concepts are nodes linked by labeled arcs. Workers in arti-

ficial intelligence had observed that a hierarchy of nodes linked by ISA relations is an efficient storage system: since all of the properties attributed to a superordinate node are inherited by its hyponyms, those properties need to be stored only once rather than separately with every hyponym. For example, when you are told that Cuthbert is a cat you know immediately that Cuthbert purrs, has four legs, fur, retractable claws, and so on. It is not necessary to learn each property separately.

There is general agreement among psychologists and workers in artificial intelligence on the idea that some kind of semantic hierarchy is required in order to represent lexical knowledge.

Constructing the hyponymic hierarchy for a broad sample of the English lexicon is a formidable task. Much of the information required is contained in the defining phrases of standard dictionaries, where a common form of definition is: "x is a y that P," where x is a hyponym of y and P is a relative clause that distinguishes x from the other hyponyms of y. For example, THE LONGMAN DICTIONARY OF CONTEMPORARY ENGLISH says that a TREE is "a type of tall PLANT with a wooden trunk and branches, that lives for many years," from which it is obvious that TREE is a hyponym of PLANT.

This kind of information can be extracted from a machine-readable dictionary (Amsler 1980, 1981; Amsler & White 1979; Chodorow, Byrd & Heidorn, 1985). The results make it clear that lexicographers work with a fundamentally consistent semantic hierarchy. Unfortunately, definitions in standard dictionaries are not written with this analysis in mind.

One feature of dictionaries that deserves comment is that it is much easier to identify superordinates from the defining phrase than to identify hyponyms. For example, the definition of *tree* will almost always say that a tree is a *plant*, but it will not go on to say that *apple*, *elm*, *fir*, *maple*, etc. are all trees; for that information a user must consult the individual entries to *apple*, *elm*, etc., which presupposes he already has the information he is searching for. In WordNet, moving down the hyponymic hierarchy should be as easy as moving up.

The hyponymic hierarchy is also apparent in standard thesauruses: ROGET'S INTERNATIONAL THESAURUS has 6 to 8 tiers of categories, going progressively from highly abstract generic categories to highly concrete specific categories. However, Roget and his successors were not slavishly devoted to the hyponymic relation, and careful judgment is sometimes required in order to extract the hyponymic relation from all the other information in an entry. Sedelow and Sedelow (1986) comment that there is much greater descriptive and analytic power, semantically, in the lower tiers of ROGET'S THESAURUS.

In most cases, the judgment required to settle questions about hyponymic relations are not difficult. In order to decide whether x is a hyponym of y, substitute them into a standard frame of the form: x *ISA* y, then judge whether the resulting proposition is true or false. If it is true, then x can be accepted as a hyponym of y.

By using a collection of dictionaries and thesauruses, liberally seasoned with linguistic intuitions, WordNet editors have introduced hyponymic relations into the synonym sets with relatively little trouble. In some cases, a word that seems to have no obvious synonym can be tied into the semantic structure through its superordinate. *Blunderbuss*, for example, has no good synonym in English, but it can be integrated into WordNet as a hyponym of *firearm*. In other instances, an initial synonym set can be reorganized: coordinate terms—names of trees, for

example—that were entered initially as a synonym set could, with the introduction of hyponymic relations, be entered more accurately as hyponyms — in this example, as hyponyms of *tree*. In general, the addition of hyponymy has had the effect of sharpening the semantic distinctions that can be drawn and, as a consequence, reducing the average size of the synonymy sets.

### *Antonymic Clusters*

Psychologists also have an interest in antonymy, since antonyms are so often used to anchor the ends of scales used in subjective judgments: *good-bad*, *agree-disagree*, *right-wrong*, etc. Probably the most extensive use of antonyms for scaling purposes was Osgood, Suci, and Tannenbaum's (1957) attempt to map all concepts into a space whose coordinates were given by pairs of antonymous adjectives.

Not every word has an antonym, of course. This relation is probably clearest between adjectives, although it is by no means limited to adjectives. The adjectival synonym sets were chosen as the most appropriate place to introduce antonymy into WordNet.

The work began with the assumption that antonymy and synonymy are themselves opposites. That is to say, synonyms are words whose meanings are very similar, whereas antonyms are words whose meanings are very dissimilar. That assumption may suffice as long as one does not look too closely, but careful analysis reveals important differences. The long history of disagreement about the nature and definition of antonymy (Egan, 1984) should have been a warning, but the extent of the difference was not recognized until an attempt was made to represent antonymous pairs by symmetrical cross-references between contrasting synonym sets.

The design of WordNet landed it, inadvertently, in the middle of a traditional argument about antonymy. Is an antonym (1) any one of several words that can be opposed to a group of synonymous terms, or is it (2) a single word, or at most one of two or three words, that can be opposed to a given word? As originally conceived, WordNet incorporated assumption (1). That is to say, relatively large groups of synonyms were first compiled; then attempts were made to cross-reference the antonymous sets. But it proved difficult to carry that program through. When synonym set  $C_i$  was put in opposition to synonym set  $C_j$ , not every word in  $C_i$  was an antonym of every word in  $C_j$ , and vice versa, and that fact made it difficult to judge whether the concepts represented by the synonym sets were truly antonymous.

For example, the concept that is represented by the synonym set { *damp*, *dank*, *drenched*, *moist*, *soaked*, *waterlogged*, *wet* } seems to be antonymous to the concept that is represented by the synonym set { *arid*, *baked*, *dehydrated*, *desiccated*, *dry*, *parched*, *sere*, *withered* }, but few people would think of *withered* as an antonym of *waterlogged*, say, or, of *baked* as an antonym of *dank*, etc. Assumption (1) defines antonymy as a relation between lexical concepts, whereas assumption (2) defines antonymy as a relation between words. Judgments of antonymy are much easier to make between words than between concepts.

The addition of antonymous relations sharpens considerably the semantic distinctions that are required. That is to say, the adoption of assumption (2) necessar-

ily limits the number of words in many synonym sets to two or three. But the notion that antonymy is a relation between words, rather than between concepts, finds support in the frequent use of derivational morphology to signal antonymy: *perfect-imperfect*, *advantageous-disadvantageous*, *benevolent-malevolent*, *powerful-powerless*, *superior-inferior*, *definite-indefinite*, etc. Or, to put it differently, prefixing *un-* to adjectives can result in new adjectives (*pleasant-unpleasant*) in much the same way that adding *en-* to adjectives can result in causative verbs (*rich-enrich*). In both cases the affix does important semantic work, but both dyads reflect formal relations between pairs of words. This is consistent with assumption (2), which defines antonymy as a relation between words.

Moreover, if it is assumed that the morphological relations involved in particular antonymous pairs must be learned by repeated exposure and practice, much the way all formal (i.e. phonological and morphological) features of English are learned, then other observations about antonyms could be explained. For example, although *big-little* and *large-small* are both antonymous pairs, it sounds odd to cross them: *big-small* and *large-little*. The explanation is that we have heard them paired one way much more frequently than the other. Although the cross is conceptually correct, it is morphologically unfamiliar.

How can a conceptual definition of synonymy coexist with a formal conception of antonymy? Or, in more practical terms, how can a loose definition of synonymy be combined with a strict definition of antonymy? Solving this practical problem, forced an interesting structure onto the adjective file: antonym pairs must form the basic skeleton of adjectival semantics, and this skeleton is fleshed out by those adjectives that have no obvious antonym but are similar to adjectives that do have antonyms. That is to say, another relation, dubbed semantic similarity, is introduced to preserve sets of several synonyms, but without precluding the one:one pairing of antonyms.

The result is illustrated in Table 1 by the cluster of concepts around the antonymous pair *wet-dry*. (The 'a' following each number indicates that it is the name for an adjectival synonym set.) If *dry* in 1005a is consulted in search of an antonym, *wet* will be found in 1000a (and vice versa), whereas if *dry* in 1015a is consulted, the antonym in 1070a will be *sweet*. On the other hand, if 1005a is consulted for near synonyms of *dry*, all the words in 1006a, 1007a, 1008a, 1009a, and 1014a will be found. Thus, a narrow interpretation of antonymy can coexist with a broad interpretation of synonymy. Moreover, this form of representation poses no special problems for polysemous words: the *dry* that is the antonym of

Table 1

The antonymic cluster, *wet-dry*  
(Antonymic relation, \*; similarity relation, &)

1000a	{ <i>wet</i> , &1001a, &1002a, &1003a, *1005a, }
1001a	{ <i>damp</i> , <i>dank</i> , <i>moist</i> , &1000a, }
1002a	{ <i>drenched</i> , <i>saturated</i> , <i>soaked</i> , <i>waterlogged</i> , &1000a, }
1003a	{ <i>foggy</i> , <i>humid</i> , <i>misty</i> , <i>rainy</i> , &1000a, }
1004a	{ <i>drunk</i> , <i>slopped</i> , <i>tipsy</i> , <i>wet</i> , *1080a, }
1005a	{ <i>dry</i> , *1000a, &1006a, &1007a, &1008a, &1009a, &1014a, }

1006a	{ arid, &1005a, }
1007a	{ dehydrated, dessicated, sere, withered, &1005a, }
1008a	{ baked, parched, &1005a, }
1009a	{ thirsty, &1005a, }
1010a	{ dry, impassive, matter-of-fact, unemotional, *1020a, }
1011a	{ barren, dry, sterile, unproductive, *1030a, }
1012a	{ boring, dry, insipid, wearisome, &1040a, &1090a, }
1013a	{ bare, dry, plain, unadorned, *1060a, }
1014a	{ anhydrous, &1005a, }
1015a	{ dry, &1110a, *1070a }
1020a	{ emotional, *1010a, }
1030a	{ fruitful, productive, *1011a, }
1040a	{ dull, &1012a, &1090a, }
1050a	{ interesting, *1090a, }
1060a	{ adorned, fancy, *1013a, }
1070a	{ sweet, *1015a, &1100a, }
1080a	{ dry, sober, *1004a }
1090a	{ uninteresting, &1012a, &1040a, *1050a, }
1100a	{ sugary, *1110a, &1070a, }
1110a	{ sugarless, &1015a, *1100a, }

*wet* expresses a different concept from the *dry* that is the antonym of *sweet*, and different also from the *dry* that is similar to *dull* and *uninteresting*.

Implicit in the adoption of this structure for WordNet is the hypothesis that native speakers of English have a similar organization of their lexical memory for antonyms. That hypothesis was explored in a series of experiments by Gross, Fischer, and Miller (1989). The first experiment asked native speakers of English to judge relations between different types of contrasting pairs of adjectives: direct antonyms, indirect antonyms, and unrelated adjectives. Direct antonyms are lexically opposed terms such as *wet* vs. *dry*. An indirect antonymic pair consists of an adjective and a near synonym of its direct antonym that does not have its own lexical antonym: *dank* vs. *dry*. Examples of unrelated adjectives are *pleasant* vs. *scarlet* or *regretful* vs. *clumsy*. Native speakers of English were expected to judge direct antonymic pairs like *wet-dry* faster than indirect pairs like *wet-parched*. A second experiment asked subjects to distinguish direct antonyms from all other types of adjective pairs. The results of these experiments, although not as robust as might have been expected, were consistent with the hypothesis that semantic memory for adjectives is organized around bipolar attributes and that certain pairs (the direct antonyms) label the poles.

Such experiments serve to illustrate one way that WordNet contributes to our understanding of the organization of lexical memory.

### *Meronymy*

Meronymy, the part-whole relation, is another basic semantic relation between words and concepts. This relation turns out to play a prominent role in the noun component of the lexicon and is widely exploited in WordNet. Winston, Chaffin

and Herrmann (1987), also Chaffin, Herrmann and Winston (1988) studied a wide variety of part-whole relations.

The most easily identifiable examples of meronymy are found among words denoting concrete and countable entities. Body parts, for example, lend themselves well to part-whole classification: a *finger* is a part of a *hand*, a *hand* is a part of an *arm*, and an *arm* is a part of a *body*.

Another kind of meronymy is represented by those cases where the concept of the whole exists only by virtue of the existence of a multiple of the parts and is conceptually and linguistically inseparable from them, as in the example *a tree is a part of a forest*. Thus, one can say *a forest is many trees* but not, for example, *a body is many arms*.

In the lexicon of nouns referring to substances, meronymy again takes on a slightly different meaning. As Lyons (1977) points out, *gold* is a substance and it can also be a part of a compound matter. Thus, we can say both *this substance is gold* and *gold is part of this substance*. But the same does not hold for *arm*: although we can say *The finger is part of an arm*, we cannot say *This arm is a finger*.

Meronymy overlaps with hyponymy in the case of collective nouns such as *furniture*: while *table* is a kind of *furniture*, it is also part of *furniture*, in the sense that the concept *furniture* can be said to prototypically include the concept *table*. The classification of such collectives can, therefore, be problematic.

In the realm of concrete and count nouns, meronymy permits the establishment of hierarchical structures in parallel with, but distinct from, hyponymic structures. Meronymic relations, like hyponymic relations, are also transitive, in that we can say that if *x* is a part of *y*, and *y* is a part of *z*, then *x* is also part of *z*. For example, a *foot* is a meronym of *leg* and *leg* is a meronym of *body*; therefore, *foot* is a meronym of *body*. It would be interesting to test whether and how meronymic transitivity is represented in lexical memory: e.g., to see whether subjects will easily associate two words that are distantly related by meronymy such as *doorknob* and *house*, and if such associations require more time than those between less distantly related words like *door* and *house*.

Interesting relations exist between the hyponymic hierarchy and the meronymic hierarchies. For example, it is not necessary to say that *deck* is a meronym of *warship* if it has already been said that *deck* is a meronym of the superordinate *ship*. Tversky and Hemenway (1984) argue that the appropriate level in the hyponymic hierarchy for entering part-whole relations is the level that has been called "basic" by anthropological linguists (Berlin, Breedlove, & Raven 1966; Rosch, Mervis, Gray, Johnson, & Boyes-Braem 1976).

Hyponymy, antonymy, and meronymy reflect different aspects of the organization of human lexical memory and they all differ from synonymy. Consequently, the four relations must be represented differently in WordNet. Not until experience had been gained with this task, however, was the extent of their differences and interrelations appreciated. In the final section of this paper, we discuss the role of these relations in the verb lexicon, which presents a great challenge to any lexicographer.

## Semantic Relations in the Verb Lexicon

At present, over 3,000 synonym sets of verbs have been compiled. They were initially classified into fifteen groups along the lines suggested in Miller and Johnson-Laird (1976). This classification follows very general but intuitively basic semantic criteria; thus, we have verbs of possession, communication, mental state and activity, motion, contact, change, competition, consumption, bodily functions, creation, psychological verbs, existence, social activities, perception, and natural events. The semantic relations of hyponymy, antonymy, and meronymy, that serve naturally to relate nouns and adjectives turn out to be less fitting for verbs.

Superficially, verbs do not seem to be easily represented by a hyponymic taxonomy. Rather than functioning as true hyponyms of a superordinate term, clusters of verbs seem to be related to a core or genus verb via a relation that often specifies the manner in which the subordinate is related to the superordinate. Thus, rather than bearing an ISA relation, a verb's relation to its genus term is more precisely expressible by means of a formula such as *to V1 is to V2 in some manner*. For example, *to skulk is to walk in a stealthy manner*; *to sew is to make by drawing together with a needle and thread*.

Among the nouns that have been entered into WordNet so far, the hyponymic *kind of* relation tended to be inherited fairly regularly, so that a hierarchical tree could easily be constructed. However, the corresponding relation among verbs behaves quite unpredictably with respect to inheritance. For example, while *walking* is a kind of *traveling*, and *prowling* and *skulking* are kinds of *walking*, it is rather odd to state that *prowling* and *skulking* are kinds of *traveling*. On the other hand, the hyponymy relation that exists between *walking* and *traveling* and, e.g., *marching* and *walking* is inherited, so that we can say that *marching* is a kind of *traveling*. These observations suggest that both *walking* and verbs referring to kinds of *walking* have two principal semantic components, one of displacement (*traveling*) and one of manner. The relation to the hypernym *travel* is only inherited in those verbs where displacement constitutes the salient semantic component. This indicates that the inheritance relations among verbs are less straightforward than they are among nouns.

Meronymy, which was found to play a significant role as a semantic relation among nouns, is not found in the same way among verbs. Its counterpart in the verb lexicon is a hierarchy-building relation that may be called semantic inclusion. Semantic inclusion is related to the logical notion of entailment. Thus, under the literal interpretation of *dream*, (1) entails (2), because when (1) holds, then (2) also holds:

- (1) John is dreaming.
- (2) John is sleeping.

Note that this relation is different from the one discussed above: *dreaming* is not a kind of *sleeping*. While the kind-of relation always includes the entailment relation (i.e., you cannot do *x* in a certain way without doing *x* itself), the converse is not true, as the examples above show; hyponymy and inclusion are distinct and asymmetric relations. Thus, both the ISA (kind of) and the inclusion relations build hierarchical, partially overlapping structures in the verb lexicon.

Hyponymy in the verb lexicon can be thought about in two different ways: 1) in terms of functions; and 2) in terms of taxonomies.

Using functional relations such as *in the manner of*, *by means of*, etc. yields a non-hyponymic organization that relates a set of verbs to a base verb by mapping the latter into a set of functions, yielding more specific verbs. In the following examples the base word is related to its more specified relatives by the functions *by means of*, *with a substance*, and *in the manner of*, respectively.

*attach* → {*nail, tape, paste, glue, pin, sew, button, hook . . .*}  
*cover* → {*butter, grease, mask, clothe, paint, plaster, . . .*}  
*cook* → {*broil, fry, boil, roast, bake, steam, . . .*}

Notice that the more specified verbs listed in the previous examples are not merely types of the base verb. The function that relates them to the base verb is additive.

The taxonomic perspective recognizes hyponymic relations among sets of verbs, analogous with the superordinate, basic, and subordinate levels that have been identified in the noun lexicon. A true superordinate of a basic level verb must contain the information expressed by the relational functions. Thus the superordinate of *to nail* is a verb meaning *to attach by means of some instrument*; the superordinate of *to butter* would mean *to cover with some substance*; and the superordinate of *to broil* would mean *to cook in some manner*. In contrast with noun taxonomies in English, the architecture of verb taxonomies is confounded by lexical gaps at the superordinate level.

Consider the taxonomic organization of two standardly recognized verb classes: the CREATION class and the CHANGE-OF-STATE class. Each class has a genus term (*create* and *change*, respectively) that characterizes the class as a whole and picks out a set of sub-classes, all of which share certain semantic and syntactic properties in common. Each class contains several superordinates:

#### CHANGE-OF-STATE:

*cook* [basic: {*broil, fry, boil, bake, roast . . .*}]  
*change (integrity)* [basic: {*break, smash, crack, shatter . . .*}]  
*change (shape)* [basic: {*crush, dent, bend . . .*}]  
*change (size)* [basic: {*grow, shrink, expand . . .*}]  
*change (consistency)* [basic: {*melt, liquefy, freeze . . .*}]  
 etc.

#### CREATION:

*create (by mental act)* [basic: {*contrive, invent, fabricate . . .*}]  
*create (from raw material)* [basic: {*sew, bake, knit, paint . . .*}]  
*create (by removal)* [basic: {*dig, drill, tear, bite, break . . .*}]

The broad classification of these and other verbs can moreover be motivated by certain kinds of syntactic behavior that members of each class share in common. For example, the basic level verbs of *create from raw material* easily permit definite object deletion, while those in the *create by mental act* and the *create by removal* generally do not:

*John is knitting an afghan./John is knitting.  
Sam is painting a portrait./Sam is painting.*

*Dick is fabricating lies./<sup>\*</sup>Dick is fabricating.  
Mary tore a hole in her coat./<sup>\*</sup>Mary tore.*

Many more CREATION and some CHANGE-OF-STATE verbs allow definite object deletion, though in many cases a strong context is required:

- 1a. Elaine is roasting a goose. / <sup>\*</sup>Elaine is roasting. [where Elaine is the agent]
- b. What does Jane do in the kitchen? Well, today she is roasting.
- 2a. The man invented a new kind of mousetrap. / <sup>\*</sup>The man is inventing.
- b. Thomas Edison invented the phonograph. / Don't bother Edison, he's inventing.

Object deletion is not an inherent property of these verb classes, but rather is linked to the fact that members of the superordinate level of both these classes can function as activity verbs (like *eat, read, dance, clean*). In this activity verb realization the indefinite object can be omitted (Vendler 1967; Dowty 1979; Mittwoch 1982). Notice that at the subordinate level neither class allows deletion. Object deletion is also ruled out for the superordinates, with the exception of *cook*. This superordinate is lexicalized, unlike the others in the two classes we have considered.

Allowing object deletion may be a property characteristic of a number of classes which are able to extend to become members of the activity class and to share syntactic and aspectual properties of that class. Nevertheless, the ease with which a verb can extend to another verb class may be related to its most basic class membership.

Representing verb relations with taxonomies versus functions is a matter of choice. As far as we can determine the two descriptions are equivalent. However, the basic philosophy behind the architecture of WordNet, which requires that words be represented in terms of their relation to other words and not by definitions, functions, features, or frames, favors the use of taxonomies. A strict use of taxonomies will, in certain instances, reveal lexical holes, positions in the network where a hypernym is suggested by the overall structure but lacks a lexical instantiation.

Exploring the feasibility of using taxonomies to represent the organization in the verb lexicon is the kind of experiment WordNet encourages. The idea is to see if certain models or organizations suggested by small fragments of the lexicon remain viable when we consider the lexicon as a whole. With respect to this investigation, we have not yet coded enough of the verb lexicon to know for sure whether hyponymy relations truly deserve a place in the verb lexicon.

So far as antonymy in the verb lexicon is concerned, we strongly suspect that it is, for the most part, a secondary semantic relation derived from adjectives (of manner, degree, or intensity) or from spatial relations, among which it is a primary relation. Whenever an antonymic relation cannot be imported from elsewhere in the lexicon, we might expect a verb pair to lack an antonymic relation.

Members of two verb synonym sets are antonyms if the manner relation by which they differ is antonymic. For example, *nibble* and *gorge* are antonyms because they are related to *eat* by *little, slow* and by *much, fast*, respectively.

Antonymy also shows up systematically among verbs denoting a change from one state to another where each state can be related to a quality (e.g., *lighten* and *darken* are antonyms by virtue of the antonymic relation that holds between the two adjectives from which they are derived).

The construction of WordNet is based entirely on the conceptual relations that exist between the members of the three major categories noun, verb, and adjective. A psychological model of lexical memory however should also account for relations between words belonging to different categories. Therefore, besides an antonymy relation that is imported from other lexical categories, a unidirectional inclusion relation, and verb class hyponymy, WordNet recognizes an additional linking that assigns verbs to a particular semantic domain. For example, the verb *fleece* is linked to the noun *sheep* and a polysemous verb like *beat* is readily disambiguated when associated with different semantic domains: culinary, musical, contact, competition, and so on.

A speaker's lexical knowledge is not limited to the conceptual relations described above; it also includes knowledge of the word's syntactic use. Verbs are special. They form the backbone of the sentence and link in interesting ways to members of all the categories discussed (nouns, adjectives, other verbs) as well as to function words and whole clauses. An important part of the lexical entry of a verb is its argument structure and selectional restrictions, i.e. the number and kind of nouns it occurs with in a sentence. In this special instance we stray from the basic constraint that words be represented only in relation to other words and include syntactic and semantic frames.

To include this kind of information in WordNet, each verb synonym set is matched with a frame specifying the semantic/syntactic restrictions (a combination of subcategorizations and selectional restrictions) of its members. Since WordNet is intended for use by linguistically unsophisticated users, the codings must be simple and straightforward, drawing upon lexical knowledge the user already possesses. The coding task presents some interesting theoretical challenges. It is not clear at this point how many frames will be needed to account for all the verbs on file, but it seems desirable to keep the number small by giving only generic specifications: for example, NP<sub>human</sub> V NP<sub>non-human</sub>. On the other hand, we hope that the frames and their relations to the synonym sets can be connected in some nonrandom fashion to the semantic relations among the verbs. Some of the semantic distinctions made in the relational structures of possession verbs, for example, are reflected in a systematic way. The verbs relating to HAVE<sub>poss</sub> occur in the frame NP<sub>human</sub> V NP<sub>non-human</sub> (*John owns a car.*). The subordinates of *take* and *give* are additionally specified for a prepositional phrase with NP<sub>human</sub> and *from* and *to*, respectively. Moreover, the frames show the difference between those *give* subordinates that systematically participate in the dative alternation and those that do not (NP V NP NP vs. NP V NP *to* NP).

## Conclusion

Significant semantic differences exist between the three major syntactic categories (noun, adjective, and verb). Words from the three categories enter into synonymy relations with other words, yet each category is strongly linked to one additional

predominant relation and tends to resist systematic organization by means of other relations. Furthermore, existing models based on fragments of the lexicon need to be and are being examined on a large scale. The more we explore the potential implementation of these models the more rich and complex the problem becomes. Being forced not only to consider those words that fit our theories and therefore easily come to mind, but also to account for the remainder of the lexicon as well, has yielded many insights we would never have stumbled upon otherwise. Our hope is that when considered as a whole, the massive contents of the lexicon will conspire to narrow down the possible descriptions to a more and more general, manageable, and empirically testable set. In addition, the electronic form of a comprehensively coded lexicon should allow other researchers to extend and examine the architecture of words in new and interesting ways.

## References

### *Cited Dictionaries*

- LONGMAN DICTIONARY OF CONTEMPORARY ENGLISH. 1978. P. Procter (ed.). London: Longman Group Limited.
- ROGET'S INTERNATIONAL THESAURUS. 1977. Fourth Edition. Revised by Robert L. Chapman. New York: Harper & Row.

### *Other Literature*

- Amsler, R. A. 1980. *The Structure of the Merriam-Webster Pocket Dictionary*. Ph.D. dissertation, The University of Texas at Austin.
- Amsler, R. A. 1981. 'A taxonomy for English nouns and verbs' in *Proceedings, 19th Annual Meeting of the Association for Computational Linguistics*. Stanford, California. 133—138.
- Amsler, R. A. and J. S. White. 1979. *Development of a Computational Methodology for Deriving Natural Language Semantic Structures via Analysis of Machine-Readable Dictionaries*. Final Report, NSF Project MCS77-01315. Linguistics Research Center, The University of Texas at Austin.
- Anderson, J. R. 1976. *Language, Memory, and Thought*. Hillsdale, N.J.: Erlbaum.
- Atkins, B. T., J. Kegl, and B. Levin. 1988. 'Anatomy of a Verb Entry: From Linguistic Theory to Lexicographic Practice' in *International Journal of Lexicography* 1: 84—126.
- Berlin, B., D. E. Breedlove, and P. H. Raven. 1966. 'Folk taxonomies and biological classification' in *Science* 154: 273—275.
- Chaffin, R. and D. J. Herrmann. 1984. 'The similarity and diversity of semantic relations' in *Memory and Cognition* 12: 134—141.
- Chaffin, R., D. J. Herrmann, and M. E. Winston. 1988. 'An empirical taxonomy of part-whole relations: Effects of part-whole relation type on relation identification.' *Language and Cognitive Processes* 3: 17—48.
- Chaffin, R. and L. Peirce. 1987. 'A taxonomy of semantic relations for the classification of GRE analogy items and an algorithm for the generation of GRE-type analogies.' (Unpublished manuscript.)
- Chodorow, M. S., R. J. Byrd, and G. E. Heidorn. 1985. 'Extracting Semantic Hierarchies from a Large On-Line Dictionary' in *Proceedings of the ACL* 299—304.
- Collins, A. M., and M. R. Quillian. 1969. 'Retrieval time from semantic memory' in *Journal of Verbal Learning and Verbal Behavior* 8: 240—247.
- Cruse, D. A. 1986. *Lexical Semantics*. Cambridge: Cambridge University Press.
- Dowty, D. 1979. *Word Meaning and Montague Grammar*. Dordrecht: D. Reidel.

- Evens, M. W., B. E. Litowitz, J. A. Markowitz, R. N. Smith, and O. Werner. 1983. *Lexical-Semantic Relations: A Comparative Survey*. Edmonton, Canada: Linguistic Research.
- Egan, R. F. 1984. 'Survey of the history of English synonymy' in P. B. Gove (ed.). *Webster's New Dictionary of Synonyms*. Springfield, Mass.: Merriam-Webster. 5a—31a.
- Fillmore, C. 1986. 'Pragmatically Controlled Zero-Anaphora' in *Proceedings of the Berkeley Linguistics Society 12*, Berkeley, CA. 95—107.
- Gross, D., U. Fischer, and G. A. Miller. 1989. The Organization of adjectival meanings. *Journal of Language and Memory* 28: 92—106.
- Iris, M. A., B. E. Litowitz, and M. W. Evens. 1985. 'The part-whole relation in the lexicon: An investigation of semantic primitives.' (Unpublished manuscript.)
- Jackendoff, R. 1983. *Semantics and Cognition*. Cambridge, Mass.: MIT Press.
- Katz, J. J. 1972. *Semantic Theory*. New York: Harper and Row.
- Klima, E. and U. Bellugi, 1979. *Signs of Language*. Cambridge, MA: Harvard University Press.
- Levin, B. 1985. 'Lexical Semantics in Review: An Introduction' in B. Levin (ed.). *Lexical Semantics in Review*. Lexicon Project Working Papers, no.1. MIT, Cambridge, Mass., 1—62.
- Lyons, J. 1977. *Semantics*, 2 vols. Cambridge: Cambridge University Press.
- Miller, G.A. 1985. 'WordNet: A dictionary browser' in *Information in Data: Proceedings of the First Conference of the UW Centre for the New Oxford English Dictionary*. Waterloo, Canada: University of Waterloo. 25—28.
- Miller, G.A. and P.N. Johnson-Laird. 1976. *Language and Perception*. Cambridge, Mass.: Harvard University Press.
- Mittwoch, A. 1982. 'On the Difference Between *Eating* and *Eating something*: Activities versus Accomplishments' in *Linguistic Inquiry* 13: 113—122.
- Norman, D.A. and D.E. Rumelhart (eds.). 1975. *Exploration in Cognition*. San Francisco: Freeman.
- Osgood, C.E., G.J. Suci, and P.H. Tannenbaum. 1957. *The Measurement of Meaning*. Urbana: University of Illinois Press.
- Raskin, R. 1987. 'Electronic thesauri: Four ways to find the perfect word' in *PC Magazine* 6: 275—283.
- Rosch, E., C.B. Mervis, W. Gray, D. Johnson, and P. Boyes-Bream. 1976. 'Basic objects in natural categories' in *Cognitive Psychology* 8: 382—439.
- Schank, R.C. 1972. 'Conceptual dependency: A theory of natural language understanding' in *Cognitive Psychology* 3: 552—631.
- Sedelow, S.Y. and W.A. Sedelow, Jr. 1986. 'Thesaural knowledge representation' in *Advances in Lexicology: Proceedings of the Second Annual Conference of the UW Centre for the New Oxford English Dictionary*. Waterloo, Canada: University of Waterloo. 29—43.
- Sowa, J.F. 1984. *Conceptual Structures: Information Processing in Mind and Machine*. Reading, Mass.: Addison-Wesley.
- Talmy, L. 1985. 'Lexicalization Patterns: Semantic Structure in Lexical Forms' in T. Shopen (ed.). *Language Typology and Syntactic Description*, vol.3: *Grammatical Categories and the Lexicon*. Cambridge: Cambridge University Press.
- Tversky, B. and K. Hemenway. 1984. 'Objects, parts, and categories' in *Journal of Experimental Psychology: General* 2: 169-193.
- Vendler, Z. 1967. *Linguistics in Philosophy*. Ithaca, NY: Cornell University Press.
- Winston, M.E., R. Chaffin, and D. Herrmann. 1987. 'A Taxonomy of Part-Whole Relations' in *Cognitive Science* 11: 417—444.